

Learning to infer weather states using partial observations

Jie Chao^{1,2}, Baoxiang Pan², Quanliang Chen¹, Shangshang Yang^{2,3}, Jingnan Wang^{2,4}, Congyi Nai^{2,5}, Yue Zheng⁶, Xichen Li², Huiling Yuan³, Xi Chen², Bo Lu⁷, Ziniu Xiao²

¹School of Atmospheric Sciences, Chengdu University of Information Technology, Sichuan, China

²Institute of Atmospheric Physics, Chinese Academy of Science, Beijing, China

³Key Laboratory of Mesoscale Severe Weather, Ministry of Education, and School of Atmospheric Sciences, Nanjing University, Jiangsu, China

⁴College of Computer, National University of Defense Technology, Hunan, China

⁵Institute of Geographic Sciences and Natural Resources Research, Chinese Academy of Sciences, Beijing, China

⁶Clustertech LTD, Hong Kong, China

⁷National Climate Center, China Meteorological Administration, Beijing, China

Key Points:

- Deep generative model enables accurate spatial interpolation of weather variables from sparse observations.
- The model generates probabilistic weather estimates with reliable uncertainty quantification by combining learned priors and observations.
- The model quantifies the value of observations for reducing uncertainty, guiding optimal observation network design.

Corresponding author: Baoxiang Pan, panbaoxiang@lasg.iap.ac.cn

Abstract

Accurate state estimation of the high-dimensional, chaotic Earth atmosphere marks a Sisyphean task, yet is indispensable for initiating weather forecast and gauging climate variability. While much effort is devoted to assimilating observations and forecasts to infer weather state, the inherent low-dimensional statistical structure in atmospheric circulation, shaped by geophysical laws and geographic boundaries, is underutilized as informative prior for state inference, or as reference for assessing representative of existing observations and planning new ones. We realize these potential by learning climatological distribution from climate reanalysis/simulation, using deep generative model. For a case study of estimating 2 m temperature spatial patterns, the learned distribution faithfully reproduces climatology statistics. A combination of the learned climatological prior with few station observations yields strong posterior of spatial pattern estimates, which are spatially coherent, faithful and adaptive to observation constraints, and uncertainty-aware. This allows us to evaluate each observation’s value in reducing state estimation uncertainty, and guide optimal observation network design by pinpointing the most informative sites. Our study showcases how generative models can extract and utilize information produced in the chaotic evolution of climate system.

Plain Language Summary

Accurate estimation of weather conditions across a large area is crucial but challenging due to the complex and chaotic nature of the atmosphere. Traditional methods rely on combining observations with forecasts, which can be computationally expensive and sensitive to model biases. We propose a new approach called Climate Inpainting (CLIN) that learns the inherent spatial patterns of the atmosphere from climate data using machine learning techniques. CLIN can effectively combine the learned patterns with limited observations to reconstruct complete spatial maps of weather variables, such as temperature. We demonstrate that CLIN can accurately reproduce the key spatial features and variability of temperature over East Asia. Moreover, CLIN can quantify the uncertainty in the estimated weather maps and evaluate the importance of each observation site in reducing the overall uncertainty. This information can guide the optimal design of weather station networks. Our approach showcases the potential of machine learning in utilizing the rich information contained in climate data to improve weather estimation and observation planning.

1 Introduction

The state of the Earth atmosphere, which concerns a broad range of socioeconomic sectors and the overall environment, is characterized by the spatial distribution of a specific set of physical properties, including temperature, pressure, wind speed and direction, density, concentration of water of different phases, composition of aerosol, greenhouse gas, etc (Holton & Hakim, 2012). To determine the atmosphere state at 50 km grid resolution requires estimating the value for all the above-mentioned physical properties at around $\sim 10^7$ grids (Schneider et al., 2017). Doubling the resolution increases the total number of grids by a factor of 8. This high dimensionality poses a daunting challenge for monitoring the atmosphere (Ghil, 2020).

Current operational forecasting centers routinely update their atmosphere state estimates by combining multi-source observations and previous forecasts, so as to reboot weather forecast and gauge climate variability (Carrassi et al., 2018). Ground based observations offer direct meteorological measurements, yet come with limited spatial coverage and high maintenance cost. Remote sensing offers broader spatial coverage, yet is indirect and error prone, requiring careful calibration based on ground-based observations.

Deficiencies in observation render it an ill-posed task to estimate the state of the high-dimensional Earth atmosphere, calling for strong prior to achieve feasible solution. Forecasts from previous time steps are frequently applied to serve this mission, carrying information from previous step observations to the current step via a process-based model (Wang et al., 2000). As a result, the state estimation accuracy depends on an intricate interplay among model biases, background uncertainty, and observation error, which cannot be effectively disentangled or controlled (Law et al., 2015). Moreover, to provide multi-scale background information using forecasting models requires operational run of large ensemble high-resolution numerical simulations, which is prohibitively expensive and burdensome (Toth et al., 2003; Palmer, 2017).

Is there extra information source for inferring the state of the high-dimensional, chaotic Earth atmosphere? It turns out that, the inherent low-dimensional statistical structure in atmospheric circulation, shaped by the underlying geophysical laws and quasi-static geographic boundaries, can serve as an informative prior for state inference. The Earth climate system, like any other chaotic system, is an information producer: it gradually reveals the characteristic structure of its phase space at ever-finer scales (Gilpin, 2024). By identifying and parameterizing this characteristic structure, we can potentially bypass the curse of high dimensionality, and make more efficient use of limited observations for the state inference task.

Some pioneering works have explored this direction, leveraging the inherent structure of climate data to fill in missing observations and rebuild historical climate records. For instance, Kadow et al. (2020) developed a partial convolution method to reconstruct historical global temperature patterns based on partial observations and climate simulation. Kanngießer and Fiedler (2024) applied a similar methodology to restore the spatial extent of dust plumes in cloud-masked satellite images. Most of these practices consider deterministic models, which are designed for specific “reconstruction” problem configurations, yielding deterministic results regardless of whether observations can adequately constrain the estimation uncertainty. As a result, these methodologies generalize poorly to state inference tasks where the number or layout of observations change, fail to reproduce extremes or apply for scenarios where only limited observations are available.

A solution to these dilemmas is to shift from deterministic model to probabilistic model (B. Pan et al., 2021). Specifically, we prefer to build a probabilistic model that explicitly represents the inherent statistical structure of the atmosphere as revealed by climate observations or simulations. Thereafter, we hope to effectively and efficiently combine the learned climatological prior with incomplete observations, so as to obtain strong posterior of spatial pattern estimates. This problem setup poses two stringent requirements on the underlying probabilistic model. First, the model must faithfully approximate the high-dimensional climatological distribution as generated by the chaotic evolution of climate dynamics. Second, the model must enable flexible probabilistic inference, allowing us to efficiently obtain posterior atmospheric state estimates given arbitrary observational constraints.

To fulfill these requirements, we resort to generative machine learning, in particular, probabilistic diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song, Sohl-Dickstein, et al., 2020; Kingma et al., 2021). Probabilistic diffusion models learn to approximate complex, high-dimensional probability distributions in an iterative manner, achieving unprecedented fitting capacity and controlling flexibility (B. Pan et al., 2023; Nai et al., 2024). To demonstrate the idea, we consider a case example of inferring the spatial pattern of 2 m temperature based on sparse observations from operational meteorology stations. We learn probabilistic diffusion models to approximate the climatological distribution of 2 m temperature spatial patterns from climate reanalysis or simulation data. After carefully assessing the model’s ability to reproduce climatology, we develop tools to “inpaint” arbitrary observation constraints into the sample generation process, yielding probabilistic 2 m temperature spatial pattern estimates. Finally, we ap-

ply this methodology to evaluate each observation’s value in reducing state estimation uncertainty, and guide optimal observation network design by pinpointing the most informative sites.

2 Methodology

2.1 Data and problem setup

We consider the task of inferring the spatial pattern of 2 m temperature over East Asia ($15^\circ\text{N} - 45^\circ\text{N}, 95^\circ\text{E} - 125^\circ\text{E}$), using station observations covering $\sim 1\%$ grids of the considered region. To achieve this, we learn climatological distribution of 2 m temperature spatial pattern using climate reanalysis or simulation data. The reanalysis data are hourly, 0.25° 2 m temperature data from the fifth-generation global climate and weather reanalysis (ERA5) developed at European Centre for Medium-Range Weather Forecasts (Hersbach et al., 2020, ECMWF). The simulation data are 3-hourly, 0.25° 2 m temperature historical simulation from the Flexible Global Ocean-Atmosphere-Land System Model version f3-H (Bao et al., 2020, FGOALS-f3-H), which participates in the sixth phase of the Coupled Model Intercomparison Project (Eyring et al., 2016, CMIP6). The station observation data are obtained from the Chinese National Climatic Data Center (X. Pan et al., 2021).

Formally, we denote the spatial pattern of 2 m temperature for the target region as \mathbf{x} , which is a 120×120 dimensional random variable here. Our objective is to approximate the distribution of \mathbf{x} , based on large number of samples from climate reanalysis or simulation:

$$p_{\theta^*} = \arg \max_{p_{\theta}} \sum \log p_{\theta}(\mathbf{x}) \quad (1)$$

Here p_{θ} is parameterized probability density function approximator, θ^* is the optimal parameter, optimized by maximizing the overall likelihood of p_{θ} assigned to the training samples.

Given p_{θ^*} and sparse observations, we need to provide probabilistic estimates of 2 m temperature spatial patterns, i.e., $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$. Here, \odot is dot product, \mathbf{m} is observation mask, with value 1/0 denoting the existence/absence of observations for each geogrid. $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$ should yield samples that are spatially coherent and faithful to observational constraints. Also, $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$ should offer accurate uncertainty quantification. For instance, geogrids close to observation stations should typically have low state estimate uncertainties, while distant ones have high uncertainties. Finally, we prefer $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$ to be adaptive to changes in observation configurations, such as the abortion or inclusion of observation stations, or rearrangement of station network layout. Below we illustrate how to achieve these requirements using the proposed methodology.

2.2 Learning climatology with probabilistic diffusion model

We elucidate how to learn climatological distribution of the target random variable using probabilistic diffusion model, thereafter leverage this learned prior for the inference task (Sec. 2.3). For clarity, we only cover key steps necessary for establishing our methodology. Details can be found in the literature referenced through the description.

To approximate a target distribution using probabilistic diffusion model, we train a series of deep neural networks that can be chained to establish bijective mapping between the target distribution and a prior distribution (Sohl-Dickstein et al., 2015; Ho et al., 2020). Specifically, we define the following Gaussian process:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{(1 - \beta_t)}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad (2)$$

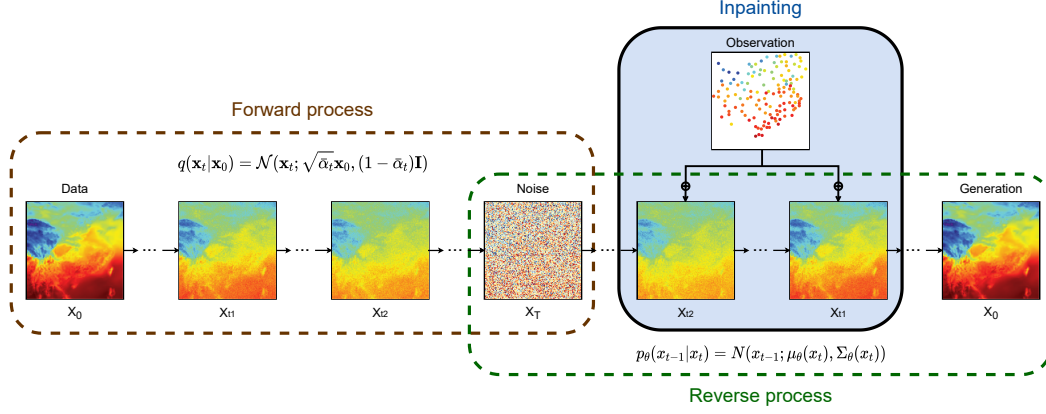


Figure 1. Overview of the Climate Inpainting (CLIN) methodology. A pre-defined forward Gaussian process (left) turns distribution of target climate variable into a prior distribution, i.e., standard Gaussian. A learned reverse Gaussian process (right) turns the prior distribution into the distribution of the target climate variable. We “inpaint” sparse observations throughout the reverse Gaussian process (right top), so as to obtain spatial pattern estimates of the target variable.

Here $p(\mathbf{x}_0) = p(\mathbf{x})$, which is the target distribution; $p(\mathbf{x}_T)$ is the prior distribution; we bridge \mathbf{x}_0 and \mathbf{x}_T using $\mathbf{x}_{t \in [1, T]}$, which are latent variables with increasing noise level; \mathcal{N} is Gaussian distribution; \mathbf{I} is identity matrix; β_t is diffusion coefficient, which is pre-defined so that, give large enough T , $p(\mathbf{x}_T|\mathbf{x}_0)$ is drawn close to $p(\mathbf{x}_T)$, which is \mathbf{x}_0 agnostic. This setup offers analytical solution for $p(\mathbf{x}_{t+\tau}|\mathbf{x}_t), \forall \tau \in [0, T-t], t \in [0, T]$, facilitating convenient inference as detailed in Sec. 2.3.

To achieve generative modeling, we reverse Eq. 2 using the following variation distributions:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta, \Sigma_\theta) \quad (3)$$

Here Σ_θ is represented as an interpolation between its analytical lower and upper bound (Dhariwal & Nichol, 2021); μ_θ can be optimized by maximizing the variational lower bound (ELBO) on the log-likelihood of the training samples (Sohl-Dickstein et al., 2015; Kingma et al., 2021). In practice, we represent μ_θ as function of neural network parameterization for $\nabla p(\mathbf{x}_t|\mathbf{x}_0)$, which is known as the *score function* (Song, Garg, et al., 2020; Song, Sohl-Dickstein, et al., 2020). This simplifies the ELBO objective function to the following form:

$$L = \mathbb{E}_{t \in [1, T], \mathbf{x}_0 \sim p(\mathbf{x}_0)} \|\nabla p(\mathbf{x}_t|\mathbf{x}_0) - \epsilon_\theta\|^2 \quad (4)$$

Here ϵ_θ is a neural network parameterization for $\nabla p(\mathbf{x}_t|\mathbf{x}_0)$. Given the trained score estimates, we can derive $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta, \Sigma_\theta)$ and sample it, starting with $p(\mathbf{x}_T)$, ending with $p(\mathbf{x}_0)$.

2.3 CLIN: inferring weather states using partial observations

We combine the learned climatology prior with station observations to infer the posterior probability distribution of the target variable, using a *repainting* methodology (Lugmayr et al., 2022; Zhang et al., 2023). Specifically, given a pre-trained diffusion model that sequentially applies $p_{\theta^*}(\mathbf{x}_t|\mathbf{x}_{t+1}) = \mathcal{N}(\mathbf{x}_t; \mu_{\theta^*}, \Sigma_{\theta^*})$ to transform $p(\mathbf{x}_T)$ to $p(\mathbf{x}_0)$, within a pre-selected time window of Ω , for grid points where we have observations, we replace

values of \mathbf{x}_t with observations noisified to time step t , by sampling $p(\mathbf{x}_t \odot \mathbf{m} | \mathbf{x}_0 \odot \mathbf{m})$. This replacement does not consider the generated parts of \mathbf{x}_t , therefore, the observations could not explicitly constrain the variability of unobserved parts.

To address this issue, for any $t \in \Omega$, after the replacement, instead of progressing to $t - 1$ directly, we rewind to time step $t - \tau$ by sampling $p(\mathbf{x}_{t-\tau} | \mathbf{x}_t)$. We thereafter repeat the denoising steps from $t - \tau$ to t for k rounds, and carry out observation replacement for \mathbf{x}_t at each round. This allows us to jointly modify both observed and unobserved regions throughout the denoising steps, yielding generated samples that are spatially coherent, faithful and adaptive to observation constraints, and uncertainty-aware. This methodology is referred to as *inpainting*, we hence name our methodology as CLIN, short for Climate Inpainting. A formal algorithm description is given below. Details for data processing, neural network architecture, hyperparameters for training and inference, are given in Supporting Information.

Algorithm 1 CLIN

Require: trained diffusion model p_{θ^*} , observations $\mathbf{x}_0 \odot \mathbf{m}$, repainting time step set Ω , rewinding step τ , rewinding round K

Ensure: observation constrained, spatially coherent sample \mathbf{x}_0

```

1: Initialize  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T - 1, \dots, 1$  do
3:    $\mathbf{x}_t \sim p_{\theta^*}(\mathbf{x}_t | \mathbf{x}_{t+1})$  ▷ Reverse sampling
4:   if  $t \in \Omega$  then:
5:     for  $k = 1, \dots, K$  do
6:        $\mathbf{x}_t^{\text{obser}} \sim p(\mathbf{x}_t \odot \mathbf{m} | \mathbf{x}_0 \odot \mathbf{m})$ 
7:        $\mathbf{x}_t \leftarrow \mathbf{x}_t \odot (\mathbf{I} - \mathbf{m}) + \mathbf{x}_t^{\text{obser}}$  ▷ Condition on observations
8:        $\mathbf{x}_{t+\tau} \sim p(\mathbf{x}_{t+\tau} | \mathbf{x}_t)$  ▷ Rewind in time by  $\tau$  steps
9:       for  $i = t + \tau - 1, \dots, t$  do
10:         $\mathbf{x}_i \sim p_{\theta^*}(\mathbf{x}_i | \mathbf{x}_{i+1})$  ▷ Reverse sampling within a rewinding round
11:      end for
12:    end for
13:  end if
14: end for
15: return  $\mathbf{x}_0$ 

```

3 Results

The accuracy for state estimation depends on 1) how well we can approximate the climatological distribution, and 2) based on a learned climatological prior, how well we can combine it with limited observations to obtain probabilistic state estimates. Below we assess model's performance for these two aspects (Sec. 3.1 and 3.2). We further employ the model to quantify the extent to which observations reduce uncertainty in state estimation, offering insights for optimal observation design (Sec. 3.3).

3.1 Climatology

We compare grid-scale and field-scale statistics of 10,000 reference/generated samples to evaluate how well the probabilistic diffusion models reproduce their training data's climatology. Two models trained with climate reanalysis (ERA5) and historical climate simulation (FGOALS) data, hereafter referred to as $\text{CLIN}_{\text{ERA5}}$ and $\text{CLIN}_{\text{FGOALS}}$, are deployed and evaluated.

The grid-scale assessment considers the mean, variance, skewness, minimum, and maximum of climatological distribution at each grid (Fig. 2). These statistics from ERA5 (Fig. 2 Row 1) and FGOALS (Fig. 2 Row 3) generally agree well, due to shared constraints from geophysical laws and geographic boundaries. The key spatial patterns are the latitudinal gradient, the influence of topography (e.g., the Tibetan Plateau), and the land-sea contrast, which are most evident in the mean, minimum and maximum maps. The variance and skewness maps reveal more regional variations. A notable discrepancy is that, compared to ERA5, FGOALS tends to hold larger skewness for most of the land regions in Southern China and Philippine Island, implying a more frequent present of high 2 m temperature for these regions.

CLIN_{ERA5} (Fig. 2 Row 2) and CLIN_{FGOALS} (Fig. 2 Row 4) can well reproduce the considered statistics of their training data, achieving high spatial correlation coefficient (~ 0.99) and low root mean squared error ($\sim 0.1^\circ\text{C}$) in matching these statistics. Besides reproducing the large scale patterns, both models accurately capture high frequency local variations influenced by complex topography, such as for mountainous regions and coastal areas. Also, the climatology difference between ERA5 and FGOALS are well reproduced by the corresponding CLIN models.

We further carry out grid-wise Kolmogorov-Smirnov tests to assess whether the generated and referential samples are likely to have come from the same underlying distribution: 96/76% grid points (stippled grids in Fig. 2) within the considered region pass a 95% confidence interval test for the CLIN_{ERA5} and CLIN_{FGOALS} model. These results suggest that the CLIN model can well reproduce climatological distribution of its training data at grid scale.

We hereafter compare the referential and generated distributions using field-scale statistics. We first examine the linear spatial structure of the 2 m temperature spatial patterns using a principal component analysis (Supporting Information Fig. S2): we decompose the spatial pattern of the target random variable into a set of orthogonal modes that capture the maximum amount of variance, and compare the spatial modes (Empirical Orthogonal Functions, EOFs), as well as the variance explained by these modes. For ERA5, the first to third leading principal components explained 90/2.7/2.0% of the total variance. While for CLIN_{ERA5}, the first to third leading principal components explained 91/2.6/1.5% of the total variance, which closely matches results for the ERA5 referential data. More importantly, we obtain spatial correlation coefficient of 0.994/0.990/0.986 between the first to third EOF of ERA5 and CLIN_{ERA5}. While the spatial modes of FGOALS differs considerably with ERA5, CLIN_{FGOALS} closely matches FGOALS: the first to third leading principal components explained 83.6/5.1/2.2% or 83.9/4.9/2.1% of the total variance for FGOALS or CLIN_{FGOALS}. The spatial correlation coefficient between the first to third EOF of FGOALS and CLIN_{FGOALS} are 0.999/0.997/0.994. These results suggest that the CLIN model can well reproduce the linear spatial mode of the considered climatological distribution.

Lastly, we examine the distribution of spatial variability across different spatial scales in the referential/generated dataset: we carry out 2D Fourier transform on the referential/generated samples, and draw the radial averaged squared magnitude of the complex Fourier coefficients as function of wave numbers (Fig. 3). The radially averaged power spectrum density of the considered referential and generated data samples follow a similar power-law scaling, suggesting that the CLIN model can well reproduce the spatial variability across scales.

To sum up, the analysis of both grid-scale and field-scale statistics demonstrates that the CLIN methodology accurately reproduces the essential characteristics and patterns of the climatological distribution present in the training data. We can thereafter leverage this learned climatological prior for the state inference task.

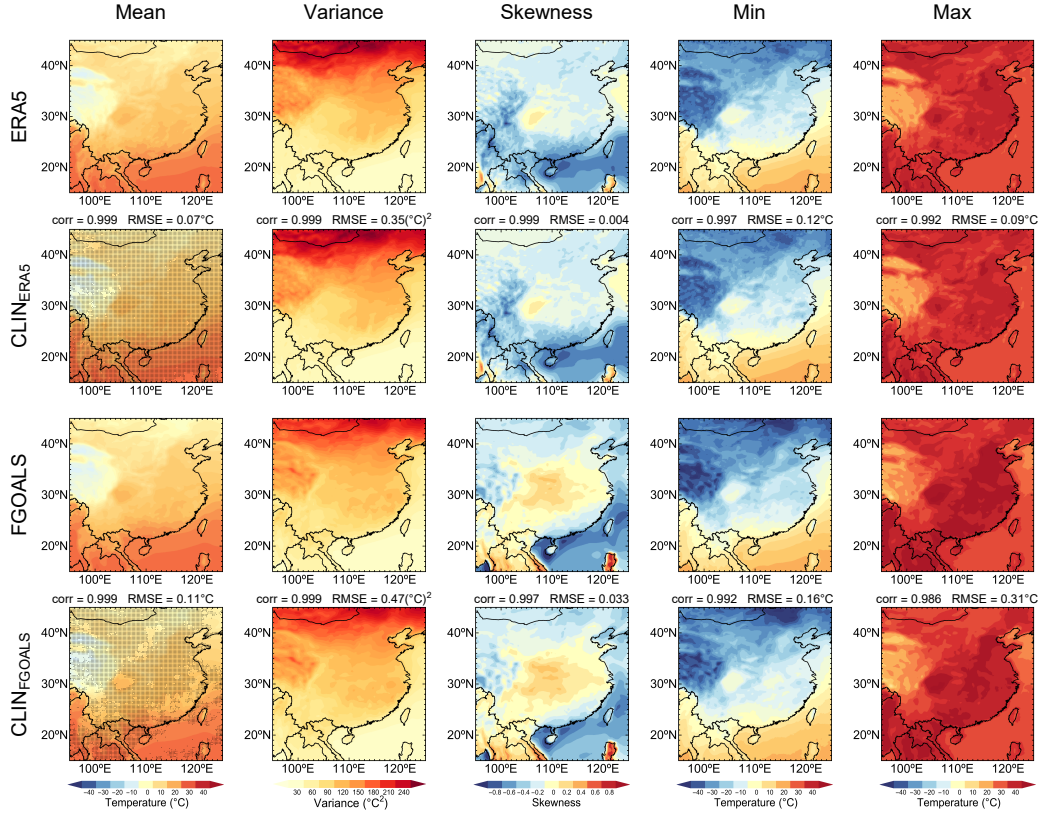


Figure 2. Grid-scale comparison of climatological statistics for climate reanalysis (ERA5, Row 1), climate simulation (FGOALS, Row 3), and probabilistic diffusion models trained using these datasets (CLIN_{ERA5}, Row 2; and CLIN_{FGOALS}, Row 4). The considered statistics are mean, variance, skewness, minimum, and maximum. The spatial correlation coefficient (corr) and root mean squared error (RMSE) between the referential dataset statistics and generated dataset statistics are labeled. Stipples denote grids that pass the Kolmogorov-Smirnov test at 95% confidence interval.

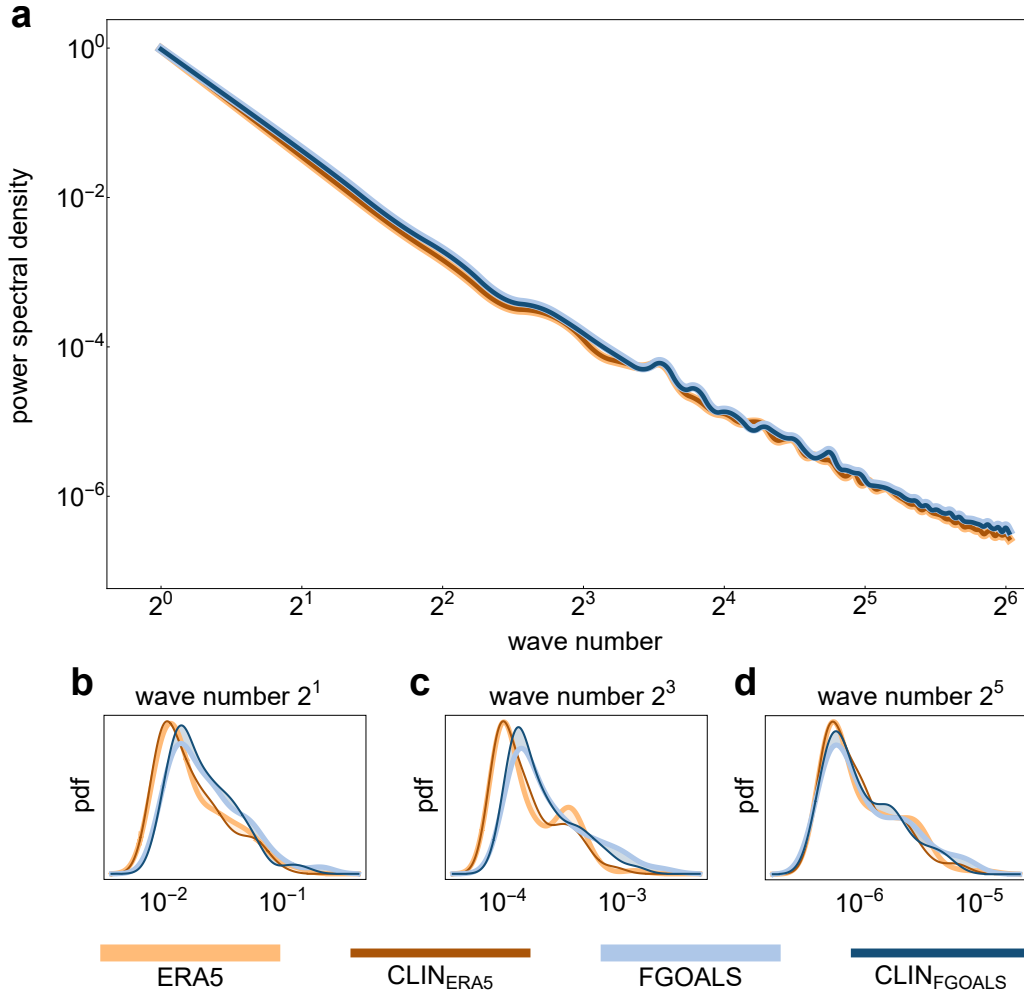


Figure 3. Radial averaged power spectrum density as function of wave number for 2 m temperature spatial pattern. **a:** results for ERA5, FGOALS, CLIN_{ERA5}, and CLIN_{FGOALS} averaged over 100 ensemble members. **b-d:** probability distribution of power spectrum density at wave number 2^1 , 2^3 , 2^5 for ERA5, FGOALS, CLIN_{ERA5}, and CLIN_{FGOALS}.

3.2 Inferring weather states using partial observations

Given a learned climatological prior, we assess how well we can combine it with partial observations to obtain probabilistic estimate of the 2 m temperature spatial patterns. The climatological priors are probabilistic diffusion models trained using climate reanalysis (ERA5) and climate simulation (FGOALS) data. The observations are from 131 operational meteorological stations across China. We randomly select 120 of these stations to inpaint into the generation process, and leave the rest 11 stations for test. For regions without station observations, we consider ERA5 data as benchmark. Below we report case example results (Sec. 3.2.1) and a 1-year round skill assessment (Sec. 3.2.2).

3.2.1 Case study

We consider four case examples covering different hours of a day and different seasons (Fig. 4). To make probabilistic inference of spatial patterns using partial observations, we gradually inpaint station observations into the generation process of $\text{CLIN}_{\text{ERA5}}$ and $\text{CLIN}_{\text{FGOALS}}$, creating 100 ensemble members for each model and each case. We report the ERA5 spatial pattern (Fig. 4 Row 1), the ensemble mean (Fig. 4 Row 2 and 5), the standard deviation of the ensemble (Fig. 4 Row 3 and 6), the mean squared error between ERA5 and the ensemble members (Fig. 4 Row 4 and 7) for $\text{CLIN}_{\text{ERA5}}$ and $\text{CLIN}_{\text{FGOALS}}$.

Both the repainted $\text{CLIN}_{\text{ERA5}}$ and $\text{CLIN}_{\text{FGOALS}}$ ensemble mean results closely match the ERA5 spatial pattern, regarding latitudinal gradient, influence of topography, and the land-sea contrast, yielding spatial correlation coefficient of $0.980 \pm 0.02 / 0.977 \pm 0.02$ for the four considered case examples. These results suggest that the proposed methodology allows effectively propagation of information from limited ($\sim 1\%$) observed locations to a broad range of unobserved parts.

Next, we test if the CLIN methodology offers reliable uncertainty quantification (Fig. 4 Row 3 and 6). A larger ensemble variance indicates greater uncertainty in the estimate, while a smaller variance suggests more confidence in the estimate. As is expected, geogrids close to observation stations tend to have low ensemble variance, while distant ones may have relatively higher ensemble variance. The information constraint from observations may be blocked by topography, such as for Tibetan Plateau and Tian Shan Mountains. While for plain regions, we can expect a larger extension of observation constraints. We further examine the relationship between the spread of the ensemble members and their estimation skill, by computing the correlation between ensemble variance and ensembles' mean squared error score. The high spread skill correlation for $\text{CLIN}_{\text{ERA5}}$ (0.90 ± 0.08) and $\text{CLIN}_{\text{FGOALS}}$ (0.94 ± 0.04) suggest that ensemble spread is a good predictor of model's estimation skill. This means that the CLIN model can capture the underlying uncertainties and provide reliable estimates of spatial estimation confidence.

To sum up, the case studies confirm that the CLIN methodology can make successful probabilistic inference of 2 m temperature spatial patterns using limited observations. The results are spatially coherent, well-constrained by observations, and offer reliable uncertainty quantification. It is worth noting that there are unneglectable mismatches between station observations and ERA5/FGOALS, regarding either climatological statistics or values. These mismatches introduce domain shift error, which is frequently encountered as we deploy a machine learning model in real-world scenarios where the data distribution differs from the training data. Below we dissect this error source by inpainting with different data sources in a 1-year round evaluation.

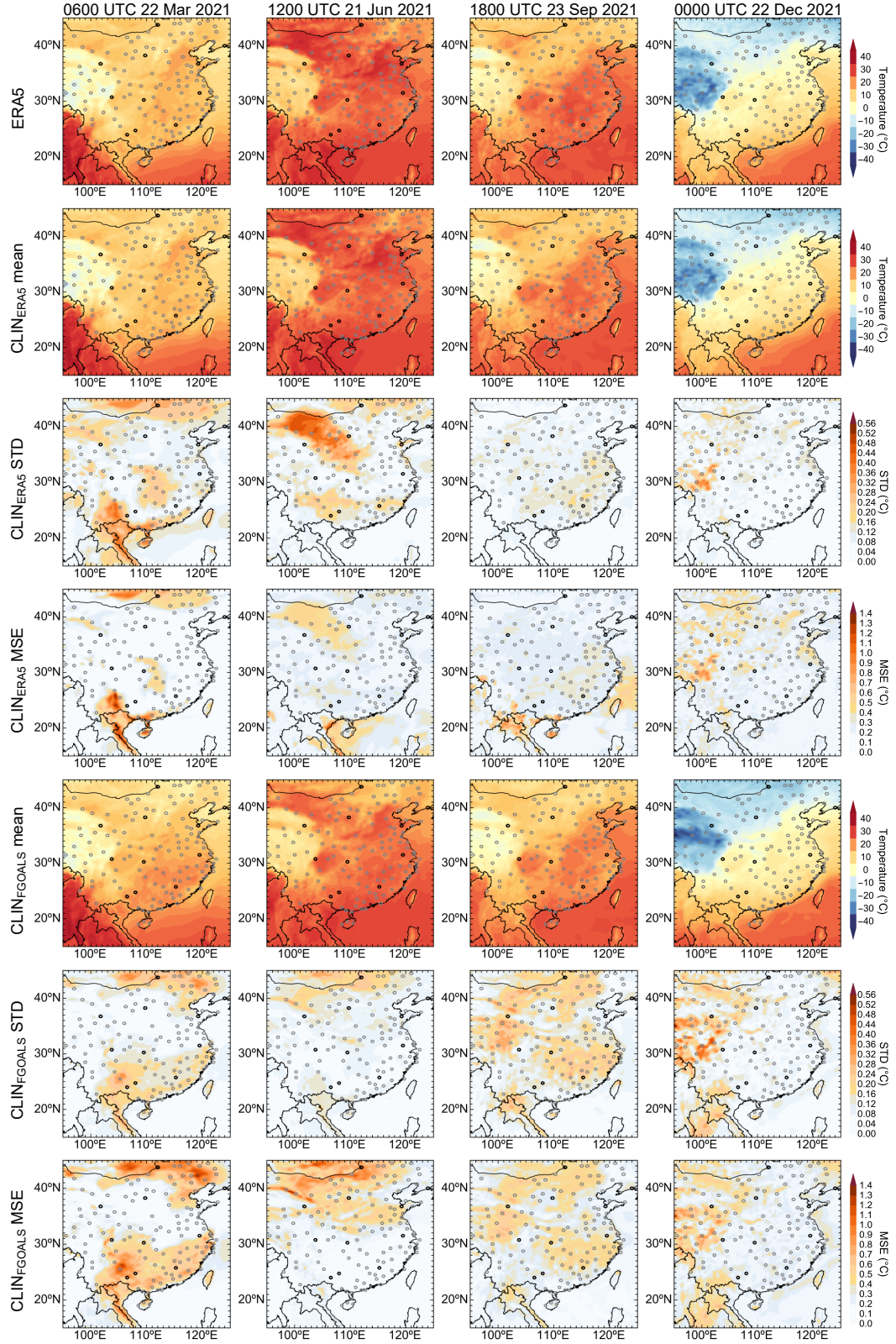


Figure 4. Case examples for probabilistic inference for 2 m temperature spatial pattern using partial observations. For $\text{CLIN}_{\text{ERA5}}$ and $\text{CLIN}_{\text{FGOALS}}$, 100 ensemble members are created by repainting observations. The ERA5 spatial pattern (Row 1), the ensemble mean (Row 2 and 5), the standard deviation of the ensemble (Row 3 and 6), the mean squared error between ERA5 and the ensemble members (Row 4 and 7) for $\text{CLIN}_{\text{ERA5}}$ and $\text{CLIN}_{\text{FGOALS}}$ are plotted.

310

3.2.2 Skill evaluation

311

312

313

314

315

316

317

We conduct a year-long evaluation of the models' performance in inferring spatial patterns, using data from Year 2021, which are not included in the models' training process. We compare ERA5 with $\text{CLIN}_{\text{ERA5}}$ and $\text{CLIN}_{\text{FGOALS}}$, both inpainted using station observations, and present the spatial distribution of their RMSE in Fig. 5a and Fig. 5b. To further investigate different uncertainty sources in the state inference task, we also consider inpainting $\text{CLIN}_{\text{ERA5}}$ using ERA5 data at the observation stations. The RMSE between this inpainted $\text{CLIN}_{\text{ERA5}}$ and the ERA5 whole-field data is shown in Fig. 5c.

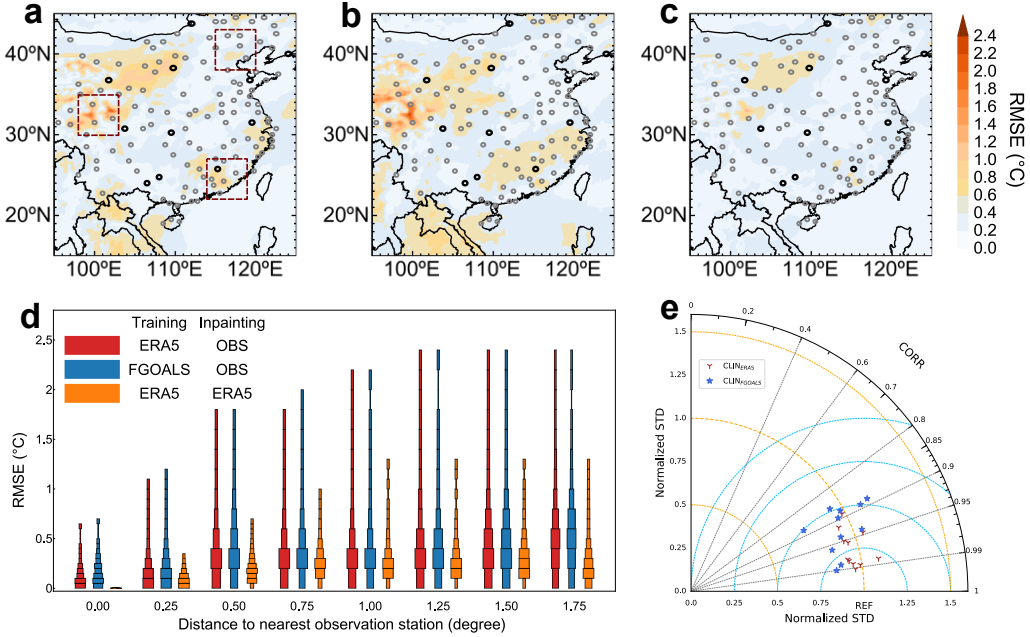


Figure 5. Skill evaluation for CLIN models to estimate spatial pattern of 2 m temperature using data for Year 2021. **a:** root mean squared error (RMSE) between ERA5 reanalysis and $\text{CLIN}_{\text{ERA5}}$ inpainted using station observations; **b:** RMSE between ERA5 reanalysis and $\text{CLIN}_{\text{FGOALS}}$ inpainted using station observations; **c:** RMSE between ERA5 reanalysis and $\text{CLIN}_{\text{ERA5}}$ inpainted using ERA5 data at station observations; **d:** distribution of RMSE as function of grid's distance to nearest observation station for the three considered methods; **e:** Taylor diagram comparing the left-out station observations with $\text{CLIN}_{\text{ERA5}}$ (orange) and $\text{CLIN}_{\text{FGOALS}}$ (blue) results. Both $\text{CLIN}_{\text{ERA5}}$ and $\text{CLIN}_{\text{FGOALS}}$ are constrained by 120 station observations here. We delineate three representative regions to evaluate the value of observations in Sec. 3.3

318

319

320

321

322

323

324

325

326

The RMSE between ERA5 and observation inpainted $\text{CLIN}_{\text{ERA5}}/\text{CLIN}_{\text{FGOALS}}$ is $0.25 \pm 0.21^\circ\text{C}/0.31 \pm 0.20^\circ\text{C}$, suggesting that the CLIN methodology enables accurate spatial pattern estimates. Both models exhibit low uncertainty in plain terrain regions or over the ocean, despite that no ocean observations were applied. This suggests that the learned climatological prior effectively captures the spatial patterns and variability in these regions, allowing the models to make confident estimates using limited and far-away observational constraints. On the other hand, both models exhibit higher uncertainty in regions with complex terrain, such as the Tibetan Plateau and the mountainous areas of Southeast China. Additionally, land areas with complicated terrain but lack-

ing observational constraints, such as Southeast Asia, also show large uncertainty in the model estimates.

The uncertainty in state inference comes from the following three sources (Tab. 1). The first is domain shift error, which is due to distribution mismatch among data applied for model training, data applied for inpainting, and data applied for skill evaluation. The second is model error, which is due to the approximation/optimization/statistical error in applying probabilistic diffusion model to fit climatological prior, or due to errors in inpainting. These two types of uncertainties are *epistemic*, as they could be reduced by gathering more data, improving the model, or incorporating knowledge about data distribution differences. The third source of uncertainty is intrinsic/aleatoric, which is due to existence of multiple plausible spatial patterns given partial observational constraints, reflecting the inherent randomness in the system being modeled.

To disentangle these uncertainty sources, we consider the following comparisons.

1. We compare the RMSE of CLIN_{ERA5} (Fig. 5a) and CLIN_{FGOALS} (Fig. 5b). CLIN_{ERA5} achieves an overall lower RMSE, which can be attributed to a relieved domain shift error from the following two aspects: a. compared to FGOALS, ERA5 better matches the “true” climatology as partially revealed by the scattered observations; b. we consider ERA5 data as “ground truth” for evaluating model performance, which gives advantage to CLIN model trained using ERA5 data.
2. We compare CLIN_{ERA5} inpainted using observation data (Fig. 5a) and CLIN_{ERA5} inpainted using scattered ERA5 data (Fig. 5c). The latter achieve significantly lower RMSE ($0.19 \pm 0.12^\circ\text{C}$), suggesting a relatively low model error and a relatively low intrinsic uncertainty of the considered task. The difference between these cases highlights the domain shift error as the observation distribution differs from ERA5.
3. We compare the performance of CLIN_{ERA5} and CLIN_{FGOALS} in predicting the observations at test stations that are excluded during repainting (Fig. 5e). For these test stations, both CLIN_{ERA5} and CLIN_{FGOALS} results show high correlation coefficient (0.87-0.99) and low root mean squared error (0.2-0.4 $^\circ\text{C}$) with the observations, with CLIN_{ERA5} performing slightly better than CLIN_{FGOALS}; CLIN_{ERA5} holds a normalized standard deviation close to 1, which closely matches the observations, while CLIN_{FGOALS} holds a normalized standard deviation slightly less than 1, suggesting a smaller temporal variability.

Table 1. Uncertainty sources for state inference using partial observations

Uncertainty source	Type	Illustration
Domain shift	Epistemic	Distribution mismatch among data applied for model training, data applied for inpainting, and data applied for skill evaluation.
Model error	Epistemic	1. Approximation/optimization/statistical error in fitting climatological prior. 2. Error in constraining the prior with observations.
Intrinsic uncertainty	Aleatoric	Existence of multiple plausible spatial patterns given observational constraints.

Finally, we quantify the spatial extension of observational constraints by showing models’ RMSE skill as function of grid’s distance to nearest observation station (Fig. 5d). We consider CLIN_{ERA5} inpainted using observation data and ERA5 data, as well as CLIN_{FGOALS} inpainted using observation data. For all these cases, models’ performances at an arbitrary grid depends closely on the grid’s proximity to observations. Meanwhile, there is large variation of models’ RMSE skills for grids that are at least 1° away from any observation stations. Below we further investigate the value of individual observations in

constraining the variability of its nearby spatial patterns, and offer guidelines for better observation planning.

3.3 On the value of observations

We apply the CLIN methodology to quantify the value of observations in constraining state estimation uncertainty, using three representative regions delineated in Fig. 5a. To achieve this, we add or remove observational stations and evaluate the impact on the estimation error (Fig. 6). Here, the first column shows the RMSE spatial pattern for the original CLIN_{ERA5} model estimates in each target region; the second column (Adding) demonstrates the impact of adding an observation station in a high-error area; the third column (Removing) illustrates the effect of removing an existing observation station; the fourth column (Terrain) provides a topographical context for each target region.

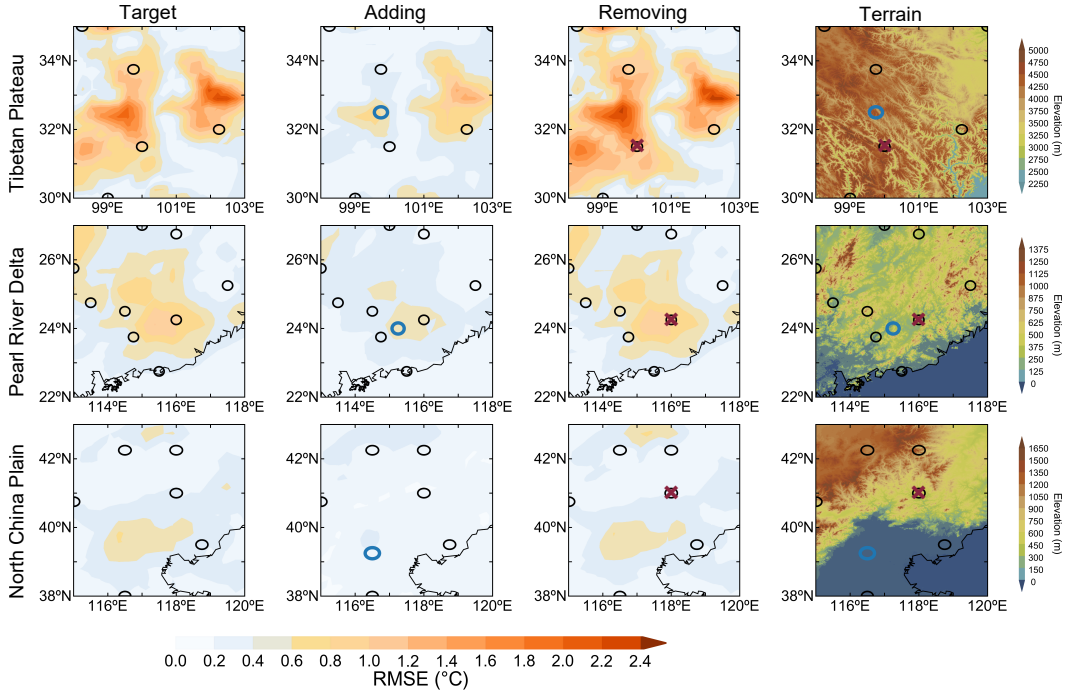


Figure 6. Evaluation of CLIN in reconstructing 2m temperature spatial pattern using different observation setups. Column 1: RMSE between CLIN_{ERA5} inpainted using observation data and ERA5 for three selected regions delineated in Fig. 5. Column 2: RMSE after including a pseudo new observation. This new observation data is from ERA5. Column 3: RMSE after including a pseudo new observation. Column 4: elevation map of the considered regions. The results are based on a year-long (Year 2021) evaluation.

For the case of Tibetan Plateau (first row), where the terrain is highly complex, with average elevations exceeding 4500 meters, we obtain a relatively high RMSE given existing observation constrains, particularly in the central and eastern parts of the region. Adding a station in the high-error area significantly reduces the RMSE for a broad range of the considered region, this impact is more pronounced here as compared to the other two cases, highlighting the importance of observational constraints in areas with complex terrain. Removing a station results in a noticeable increase in RMSE in the surrounding areas. Similarly, the effect of station removal is more evident compared to the

other two cases, suggesting that the model heavily relies on the limited observational data to constrain its estimates in this complex terrain. The loss of a station in a critical location can greatly impact the model's ability to capture the local temperature patterns.

For the case of Peal River Delta (second row), the terrain is characterized by a mix of lowlands and hilly regions, with elevations ranging from 0 to 1000 meters. The original RMSE is low overall, with some higher values in the central and northwest mountain regions. Adding a station in the high-error area effectively reduces the RMSE. Meanwhile, removing a station leads to a hardly noticeable increase in RMSE in the surrounding areas.

For the case of North China Plain (third row), the northern part is featured by mountainous terrains exceeding 1000 meters, and the southern part has flat topography and homogeneous terrain. Adding a station in the central of southern plain area reduces the RMSE significantly, as existing observations are either from the northern mountain areas, or is too far away. Same as previous case, removing a station has minimal impact on the RMSE distribution.

To sum up, we discuss the application of the CLIN methodology to evaluate the impact of observational data on state estimation uncertainty across three diverse regions. It emphasizes the importance of strategic addition and removal of observational stations in improving estimation accuracy, particularly in areas with complex terrain. The findings highlight how existing observation constraints influence RMSE distribution, with significant reductions observed when stations are added in high-error areas. Conversely, removal of stations leads to increased RMSE, underscoring the model's reliance on limited observational data. Overall, we provide valuable insights for optimizing the design of observation networks, leading to a reduction in uncertainties and biases in weather and climate analysis.

4 Conclusion

Accurate state estimation of Earth atmosphere marks a daunting task due to its high-dimensionality and chaotic nature. We demonstrated the potential of deep generative models, specifically probabilistic diffusion models, in learning the inherent low-dimensional statistical structure of atmospheric circulation from climate reanalysis and simulation data. By leveraging this learned climatological prior, we developed a methodology named CLIN (Climate Inpainting) to effectively infer weather states from partial observations.

For the case study of estimating 2 m temperature spatial patterns, the learned climatological prior accurately reproduced the essential characteristics and patterns of the training data at both grid-scale and field-scale. This learned prior effectively captured multi-scale climate patterns, providing regularization and stability to the state estimation task.

Combining the learned climatological prior with station observations, CLIN yielded strong posterior estimates of 2 m temperature spatial patterns. The estimates were spatially coherent, well-constrained by observations, and provided reliable uncertainty quantification. Regions near observation stations exhibited low ensemble variance, indicating high confidence in the estimates, while distant regions showed relatively higher ensemble variance. The high spread-skill correlation confirmed that the ensemble spread was a good predictor of the model's estimation skill.

Moreover, CLIN allowed us to quantify the value of each observation station in reducing state estimation uncertainty. By adding or removing stations and evaluating the impact on the estimation error, we demonstrated the potential of this approach in guiding the design of optimal observation networks.

Our study showcases the power of deep generative models in extracting and utilizing the information produced by the chaotic evolution of the climate system. The proposed CLIN methodology opens up new opportunities for data-driven weather state estimation, potentially complementing traditional data assimilation approaches.

Future work could focus on extending CLIN to handle indirect observations (i.e., remote sensing) and multiple interdependent variables, incorporating temporal dynamics, and adapting to long-term climate trends. Addressing the computational demands and data requirements of diffusion models is another important direction for making this approach more practical and accessible.

In conclusion, this study demonstrates the immense potential of deep generative models in advancing climate data exploration and tackling complex inference tasks in atmospheric sciences. By learning the intrinsic statistical structure of the climate system, these models can effectively bridge the gap between sparse observations and complete weather state estimates, paving the way for more accurate and efficient climate monitoring and prediction.

5 Data Availability

The ERA5 reanalysis data are obtained from the Copernicus Climate Change Service (C3S) Climate Data Store (CDS), accessible at <https://cds.climate.copernicus.eu/>.

The FGOALS model data are obtained from the Coupled Model Intercomparison Project Phase 6 (CMIP6), hosted by the Program for Climate Model Diagnosis and Intercomparison (PCMDI) at Lawrence Livermore National Laboratory (LLNL), accessible at <https://pcmdi.llnl.gov/CMIP6/>.

The observational data are freely available for download from the following website: <http://www.ncdc.noaa.gov/oa/ncdc.html>. The site information used in this study was obtained from the China Meteorological Data Network, hosted by the China National Meteorological Science Data Center (NMDC), accessible at <http://data.cma.cn/>.

6 Open Research

Model configuration, analysis scripts, data files used for this study will be publicly available upon accept of the work.

Acknowledgments

This research is supported by National Key R&D Program of China (Grant NoS. 2023YFC3007700 and 2023YFC3007705). We appreciate the insightful discussions with Dr. Niklas Boers from Technical University of Munich, and with Dr. Bin Wang and Dr. Juanjuan Liu from Chinese Academy of Science.

References

- Bao, Q., Liu, Y., Wu, G., He, B., Li, J., Wang, L., ... others (2020). Cas fgoals-f3-h and cas fgoals-f3-l outputs for the high-resolution model intercomparison project simulation of cmip6. *Atmospheric and Oceanic Science Letters*, 13(6), 576–581.
- Carrassi, A., Bocquet, M., Bertino, L., & Evensen, G. (2018). Data assimilation in the geosciences: An overview of methods, issues, and perspectives. *Wiley Interdisciplinary Reviews: Climate Change*, 9(5), e535.
- Dhariwal, P., & Nichol, A. (2021). Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34, 8780–8794.

- Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K. E. (2016). Overview of the coupled model intercomparison project phase 6 (cmip6) experimental design and organization. *Geoscientific Model Development*, 9(5), 1937–1958.
- Ghil, M. (2020). Hilbert problems for the climate sciences in the 21st century—20 years later. *Nonlinear Processes in Geophysics*, 27(3), 429–451.
- Gilpin, W. (2024). Generative learning for nonlinear dynamics. *Nature Reviews Physics*, 1–13.
- Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., ... others (2020). The era5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730), 1999–2049.
- Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33, 6840–6851.
- Holton, J. R., & Hakim, G. J. (2012). *An introduction to dynamic meteorology*. Academic press.
- Kadow, C., Hall, D. M., & Ulbrich, U. (2020). Artificial intelligence reconstructs missing climate information. *Nature Geoscience*, 13(6), 408–413.
- Kanngießer, F., & Fiedler, S. (2024). “seeing” beneath the clouds—machine-learning-based reconstruction of north african dust plumes. *AGU Advances*, 5(1), e2023AV001042.
- Kingma, D., Salimans, T., Poole, B., & Ho, J. (2021). Variational diffusion models. *Advances in neural information processing systems*, 34, 21696–21707.
- Law, K., Stuart, A., & Zygalakis, K. (2015). Data assimilation. *Cham, Switzerland: Springer*, 214, 52.
- Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., & Van Gool, L. (2022). Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 11461–11471).
- Nai, C., Pan, B., Chen, X., Tang, Q., Ni, G., Duan, Q., ... Liu, X. (2024). Reliable precipitation nowcasting using probabilistic diffusion models. *Environmental Research Letters*.
- Palmer, T. (2017). The primacy of doubt: Evolution of numerical weather prediction from determinism to probability. *Journal of Advances in Modeling Earth Systems*, 9(2), 730–734.
- Pan, B., Anderson, G. J., Goncalves, A., Lucas, D. D., Bonfils, C. J., Lee, J., ... Ma, H.-Y. (2021). Learning to correct climate projection biases. *Journal of Advances in Modeling Earth Systems*, 13(10), e2021MS002509.
- Pan, B., Wang, L.-Y., Zhang, F., Duan, Q., Li, X., Pan, X., ... others (2023). Probabilistic diffusion model for stochastic parameterization—a case example of numerical precipitation estimation. *Authorea Preprints*.
- Pan, X., Guo, X., Li, X., Niu, X., Yang, X., Feng, M., ... others (2021). National tibetan plateau data center: promoting earth system science on the third pole. *Bulletin of the American Meteorological Society*, 102(11), E2062–E2078.
- Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system modeling 2.0: A blueprint for models that learn from observations and targeted high-resolution simulations. *Geophysical Research Letters*, 44(24), 12–396.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning* (pp. 2256–2265).
- Song, Y., Garg, S., Shi, J., & Ermon, S. (2020). Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in artificial intelligence* (pp. 574–584).
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2020). Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*.

- 533 Toth, Z., Talagrand, O., Candille, G., & Zhu, Y. (2003). Probability and ensemble
534 forecasts. *Forecast verification: A practitioner's guide in atmospheric science*,
535 137, 163.
- 536 Wang, B., Zou, X., & Zhu, J. (2000). Data assimilation and its applications. *Pro-*
537 *ceedings of the National Academy of Sciences*, 97(21), 11143–11144.
- 538 Zhang, G., Ji, J., Zhang, Y., Yu, M., Jaakkola, T. S., & Chang, S. (2023). Towards
539 coherent image inpainting using denoising diffusion implicit models.